



Lecture 10: Transformer & LLM Accelerators

Notes

- Midterm grade will post tonight
- Meeting with the project teams next week

Recap

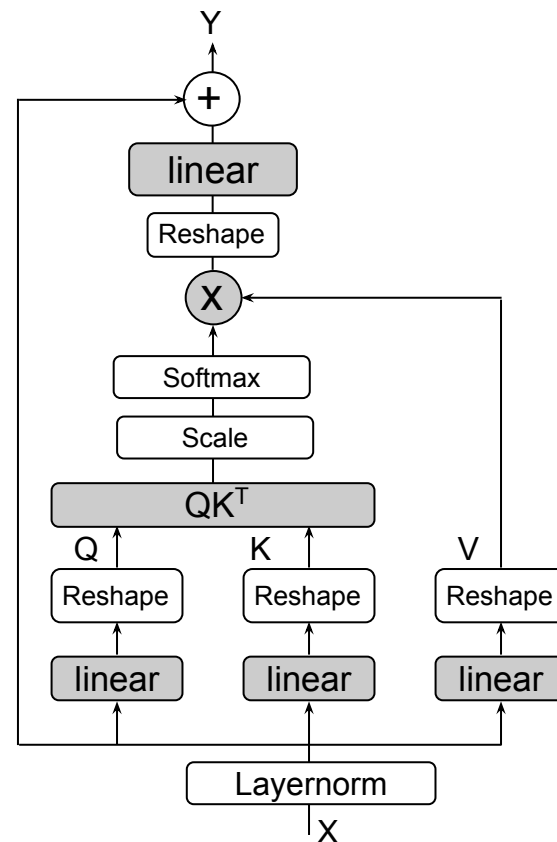
- Convolutional operation conversion
- Hardware architecture of CNN accelerator
- Systolic array
- Popular CNN accelerator design
 - SpAtten
 - EdgeBert
 - Olive

Topics

- **Matrix Multiplication with Transposition**
- Hardware design for Nonlinear Blocks
- System optimization of LLMs
- Popular transformer accelerator design
 - SpAtten
 - EdgeBert
 - Olive

Self-Attention Block

- Given input x , the first step in calculating self-attention is to create three vectors from each of the input x , denoted as: Query (Q), Key (K), Value (V).
 - $(B, L, E) * (E * E) \rightarrow (B * L * E)$
- The second step in calculating self-attention. This will compute the attention score between each pair of input tokens.
 - $QK^T \rightarrow (B, L * E) * (B, E * L) \rightarrow (B, L * L)$
- Scale and normalize the score using softmax.
 - $\text{Softmax}(QK^T) \rightarrow (B, L * L)$
- Multiply each value vector by the softmax score.
 - $\text{Softmax}(QK^T) * V$
 - $(B, L * L) * (B, L * E) \rightarrow (B, L * E)$
- Pass the result to the linear layer, sum with the input.



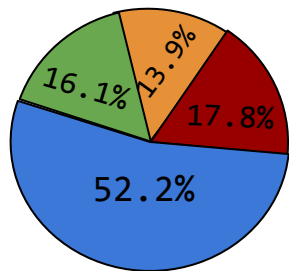
Operations Other than Multiplications

- Transposition
- Nonlinear operations
 - Softmax
 - LayerNorm
 - GeLU

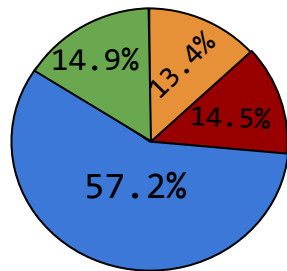
Breakdown on Computational Cost

Latency Breakdown

Matmul Normalization Softmax Others



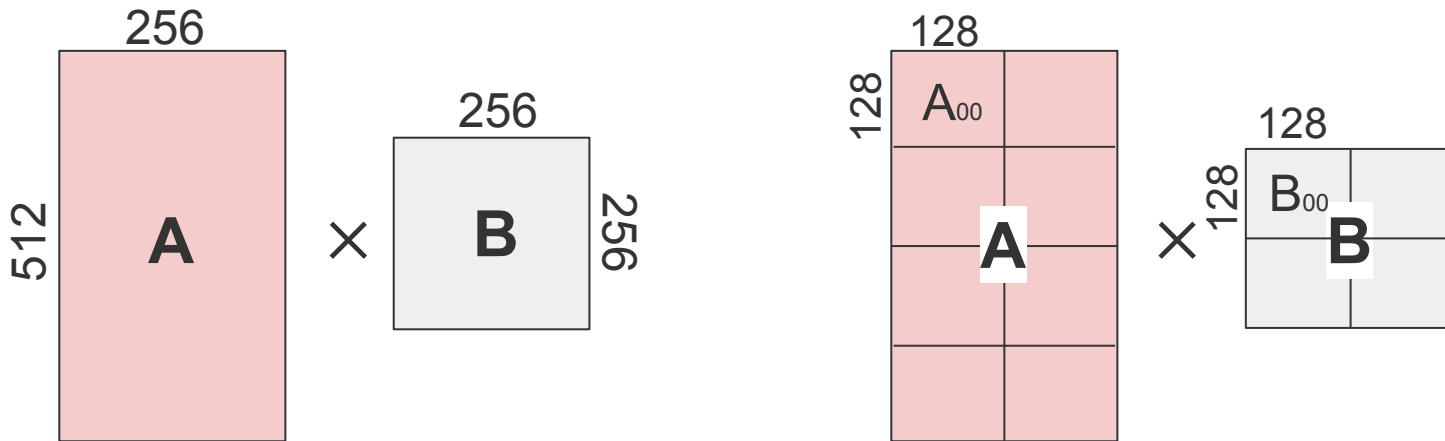
GPT2



OPT

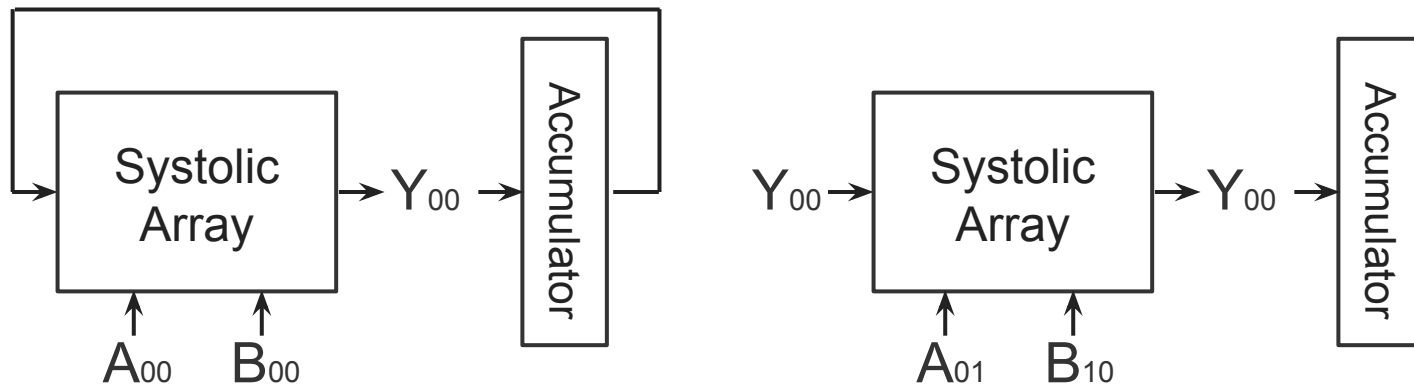
- Matmul still contributes to majority of the overall latency.
- Nonlinear operations are not negligible.
- Also other operations (e.g., transposition) also contributes to a great portion of the overall latency.

Matrix Multiplication



- The large matrix operands are first partitioned into tiles that can fit the size of the compute core.

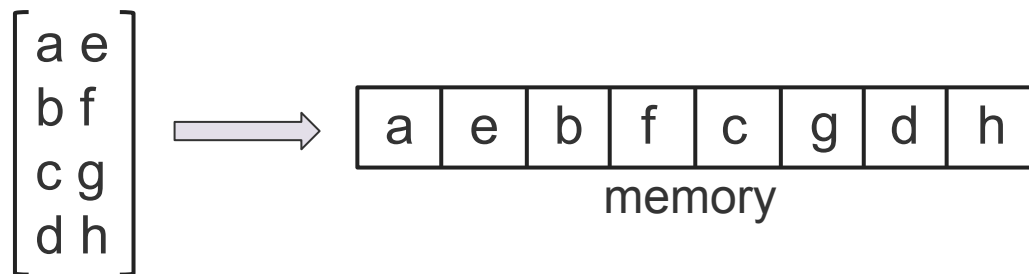
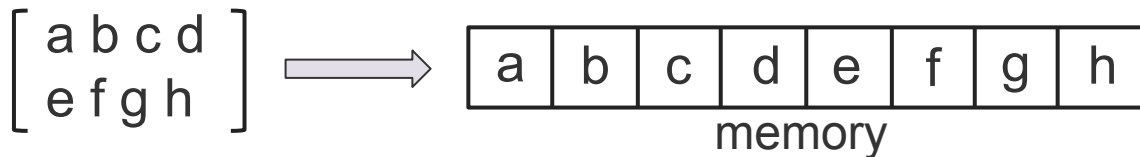
Matrix Multiplication



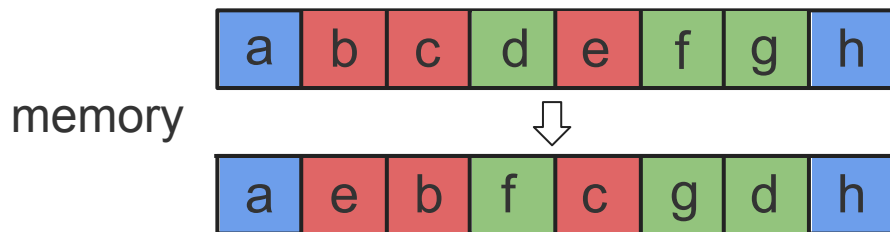
- The large matrix operands are partitioned into tiles that match the compute core's capacity, after which multiplication and accumulation are executed on a per-tile basis.
- However, sometimes the transposition operations are also required $\rightarrow QK^T$

In-Place Matrix Transposition

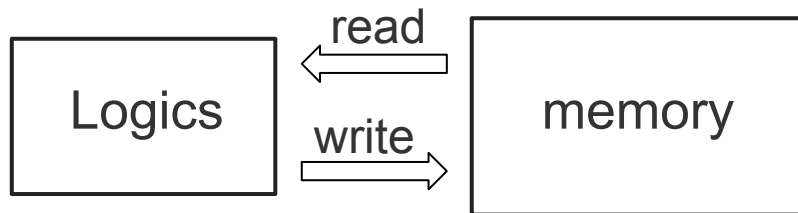
- In-place matrix transposition refers to the process of transposing a matrix directly within its existing memory space, requiring only a minimal amount of extra storage.



In-Place Matrix Transposition



$(b, c, e) \rightarrow (e, b, c)$
 $(d, f, g) \rightarrow (f, g, d)$



- Need to read multiple entries from the memory, permute them and write them back.
- This operation should be performed efficiently with minimal memory access cost.

In-Place Matrix Transposition



Step 1



Step 2

- The search for optimal swapping patterns that minimize permutations is a well-established problem in mathematics.

Topics

- Matrix Multiplication with Transposition
- **Hardware design for Nonlinear Blocks**
- System optimization of LLMs
- Popular transformer accelerator design
 - SpAtten
 - EdgeBert
 - Olive

Implementation of Nonlinear Operations: Softmax

- Softmax operations are heavily adopted in the transformer.

$$s_i = \frac{e^{z_i}}{\sum_{j=0}^{N-1} e^{z_j}} \text{ For } i = 1, 2, \dots, N$$

- For positive z with INT representation, we can approximate the values of e^z using the following derivations:

$$e^z = 2^{z \log_2 e} = 2^{u+v} \quad \log_2 e \approx 1.0111_2$$

$$z \log_2 e \approx z + (z \gg 2) + (z \gg 3) + (z \gg 4)$$

- To compute 2^{u+v} , we can perform shift and multiplication:

$$e^{z'} = 2^{u+v} \approx 2^u (1 + v/2)$$

u and v are the integer and fractional part of the exponent, $v/2$ is the mantissa, u is the exponent

Taylor Approximation

- A Taylor series is a series expansion of a function about a point. A one-dimensional Taylor series is an expansion of a real function $f(x)$ about a point $x=a$ is given by:

$$f(x) = f(a) + f'(a)(x-a) + \frac{f''(a)}{2!}(x-a)^2 + \frac{f^{(3)}(a)}{3!}(x-a)^3 + \dots + \frac{f^{(n)}(a)}{n!}(x-a)^n + \dots$$

- For small v , e^v can be approximated as:

$$e^v \approx 1 + \frac{v}{2} \quad 2^v \approx 1 + \frac{v}{2} \quad \log(1+x) \approx x$$

Taylor Approximation

- A Taylor series is a series expansion of a function about a point. A one-dimensional Taylor series is an expansion of a real function $f(x)$ about a point $x=a$ is given by:

$$f(x) = \boxed{f(a) + f'(a)(x-a)} + \frac{f''(a)}{2!}(x-a)^2 + \frac{f^{(3)}(a)}{3!}(x-a)^3 + \dots + \frac{f^{(n)}(a)}{n!}(x-a)^n + \dots$$

- For small v , e^v can be approximated as:

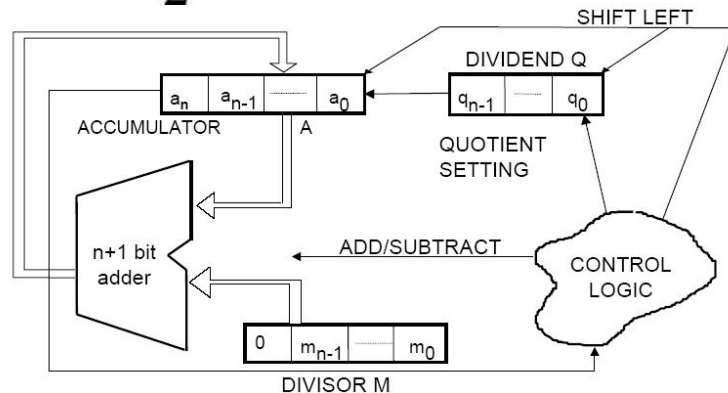
$$e^v \approx 1 + \frac{v}{2} \quad 2^v \approx 1 + \frac{v}{2} \quad \log(1+x) \approx x$$

Division

- To implement division operation with FP format, we can always apply the following derivations:

$$\frac{a}{b} = 2^{e_a} (1 + m_a) / 2^{e_b} (1 + m_b) = 2^{e_a - e_b + \log_2(1 + m_a) - \log_2(1 + m_b)}$$
$$\approx 2^{e_a - e_b + m_a - m_b} \approx 2^{e_a - e_b} \left(1 + \frac{m_a + m_b}{2}\right)$$

- For INT division, we can also implement the hardware divisor.



Implementation of Nonlinear Operations:

LayerNorm

- For the input vector z , the normalization operation requires to compute its mean and variance, then the intermediate results are scaled with some predefined values.

$$\mathbf{s} = \alpha \frac{z - \mu_z}{\sigma_z} + \beta \quad \mu_z = \frac{\sum_i z_i}{N} \quad \sigma_z = \sqrt{\frac{\sum_i (z_i - \mu_z)^2}{N}}$$

- Most of the operations are supported, the inverse of square root can be computed as follows:

$$y = \frac{1}{\sqrt{x}} \quad \log_2(y) = -\frac{1}{2} \log_2(x)$$

Implementation of Nonlinear Operations: LayerNorm

- Most of the operations are supported, the inverse of square root can be computed as follows:

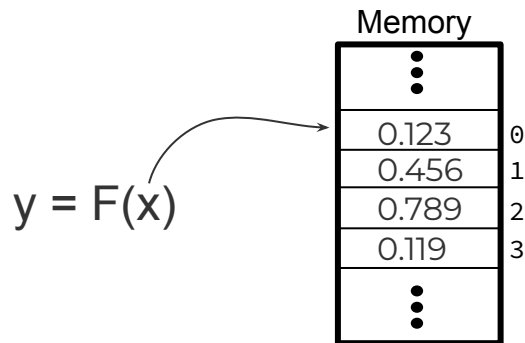
$$y = \frac{1}{\sqrt{x}} \quad \log_2(y) = -\frac{1}{2} \log_2(x)$$

$$x = 2^{E_x - Q} (1 + M_x / 2^L) \quad \log_2 x = E_x - Q + \log_2(1 + M_x / 2^L) \\ \approx E_x - Q + M_x / 2^L + \sigma_x$$

- Q is the bias, Ex and Mx are the binary representations of the exponent and mantissa, respectively.

Table Lookup

- For other complicated nonlinear functions, we can always precompute the results and store them in the buffer.



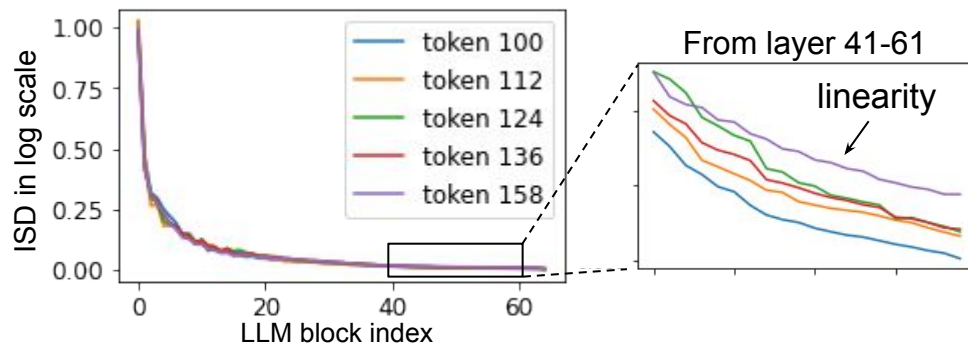
- However, this will inevitably lead to additional memory access cost and footprint.

HAAN: LayerNorm Accelerator

Layer Normalization:

$$\mathbf{s} = \alpha \frac{\mathbf{z} - \mu_z}{\sigma_z} + \beta$$

Computing the inverse of standard deviation of costly



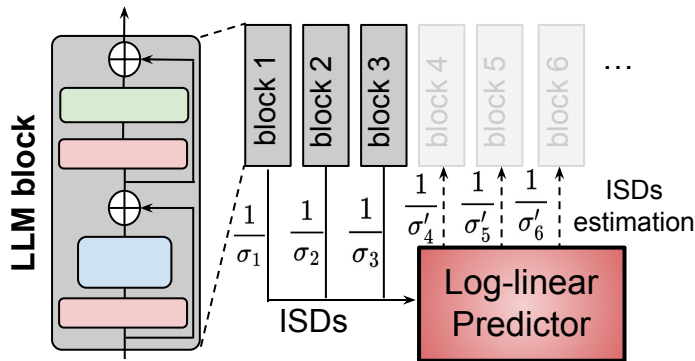
- Exploit correlation in input statistics across layers.
- Skip redundant computations and estimate normalization statistics.

HAAN: LayerNorm Accelerator

Layer Normalization:

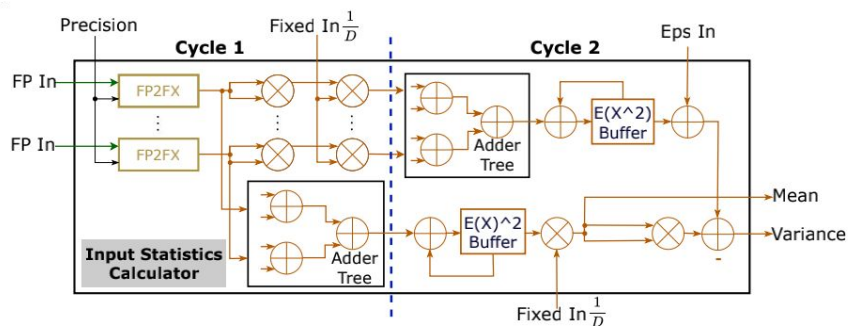
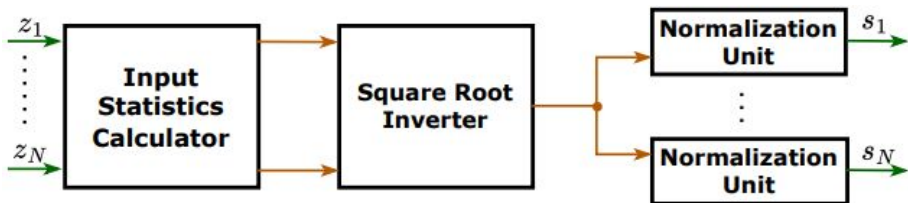
$$\mathbf{s} = \alpha \frac{\mathbf{z} - \mu_z}{\sigma_z} + \beta$$

Computing the inverse of standard deviation of costly



- Exploit correlation in input statistics across layers.
- Skip redundant computations and estimate normalization statistics.

HAAN: LayerNorm Accelerator



- **Overall Architecture**

- Input Statistics Calculator.
- Square Root Inverter.
- Normalization Unit.

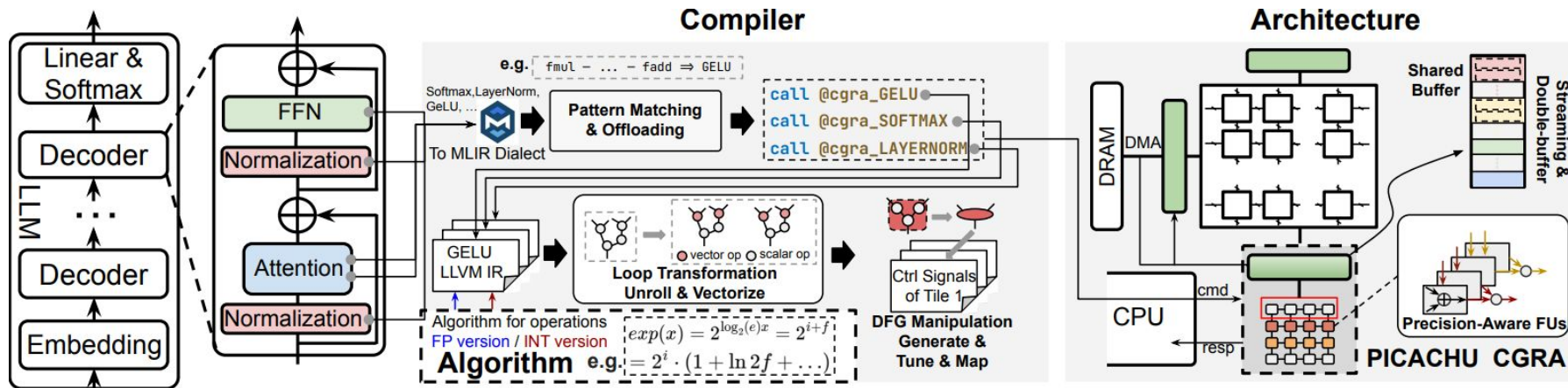
- **Input Statistics Calculator**

- Compute mean and variance.
- Parallel processing to reduce latency.

- **Square Root Inverter**

- Approximate inverse square root using Newton's method.
- Support for layer skipping.

PICACHU



- PICACHU is a plug-in coarse-grained reconfigurable accelerator tailored to efficiently handle nonlinear operations by using custom algorithms and a dedicated compiler toolchain.

PICACHU

Categories	Nonlinear Operations	Mathematical Operator	Representative LLMs
Activation Function	$\text{Softmax}(x_i) := \frac{\exp(x_i)}{\sum_{j=1}^k \exp(x_j)} = \frac{\exp(x_i - u)}{\sum_{j=1}^k \exp(x_j - u)};$ $u = \max_{j=1} x_j$	Division, Exponential	All
	$\text{ReLU}(x) := \max(0, x)$	Maximum	OPT [145], T5 [90]
	$\text{GeLU}(x) := 0.5x \left(1 + \text{Tanh}(\sqrt{2/\pi}(x + 0.044715x^3)) \right);$ $\text{Tanh}(x) = (\exp(x) + \exp(-x)) / (\exp(x) - \exp(-x))$	Division, Exponential	GPT [14, 84, 87, 88], BLOOM [57], Falcon [83], PanGu- α [144], Jurassic-1 [64], Gopher [89]
	$\text{GeGLU}(x) := \text{GeLU}(xW + b) \oplus (xV + c)$	Division, Exponential	LaMDA [110], GLM-130B [143]
Normalization Function	$\text{LayerNorm}(x_i) := \frac{x_i - \mu}{\sigma};$ $\mu = \frac{1}{C} \sum_{i=1}^C x_i, \sigma = \sqrt{\frac{1}{C} \sum_{i=1}^C (x_i - \mu)^2 + \epsilon}$	Inverted Square Root	GPT [14, 84, 87, 88], BLOOM [57], BERT [20], OPT [145], PanGu- α [144], Jurassic-1 [64]
	$\text{RMSNorm}(x_i) := \frac{x_i}{\sigma}; \sigma = \sqrt{\frac{1}{C} \sum_{i=1}^C (x_i)^2 + \epsilon}$	Inverted Square Root	LLaMA [113, 114], T5 [90], Mistral [43], Qwen [7], DeepSeek [11], Gopher [89]
Positional Embedding	$\text{RoPE} \begin{pmatrix} x_{2i-1} \\ x_{2i} \end{pmatrix} = \begin{pmatrix} x_{2i-1} \cos(m\theta_i) - x_{2i} \sin(m\theta_i) \\ x_{2i-1} \sin(m\theta_i) + x_{2i} \cos(m\theta_i) \end{pmatrix};$ $\theta_i = 10000^{-2(i-1)/d}, i \in [1, 2, \dots, d/2]$	Sine, Cosine	GPTNeo-20B [13], LLaMA [113, 114], PaLM [17], GLM-130B [143], Qwen [7], DeepSeek [11]

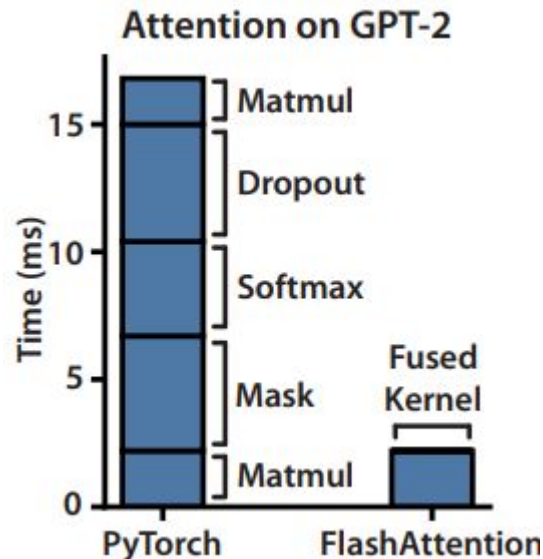
- All nonlinear operations within LLM can be broken down into various mathematical operators.

Topics

- Matrix Multiplication
- Hardware design for Nonlinear Blocks
- System optimization of LLMs
- Popular transformer accelerator design
 - SpAtten
 - EdgeBert
 - Olive

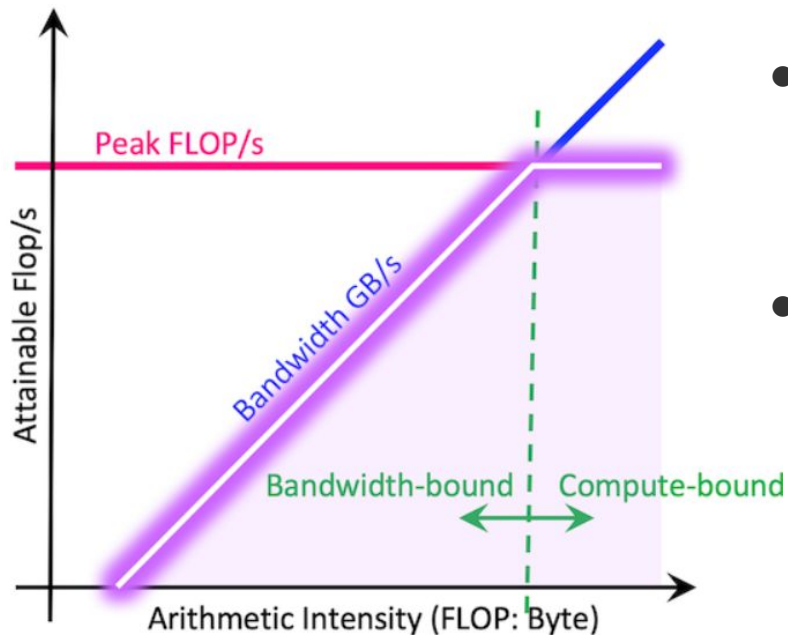
Flashattention

- Most of the operations are bottlenecked by memory speed.
- A new attention algorithm that computes exact attention with far fewer memory accesses.
- The main goal is to avoid reading and writing the attention matrix to and from memory.
- Flashattention enables to compute the softmax reduction without access to the whole input.



7.6x reduction on GPU latency

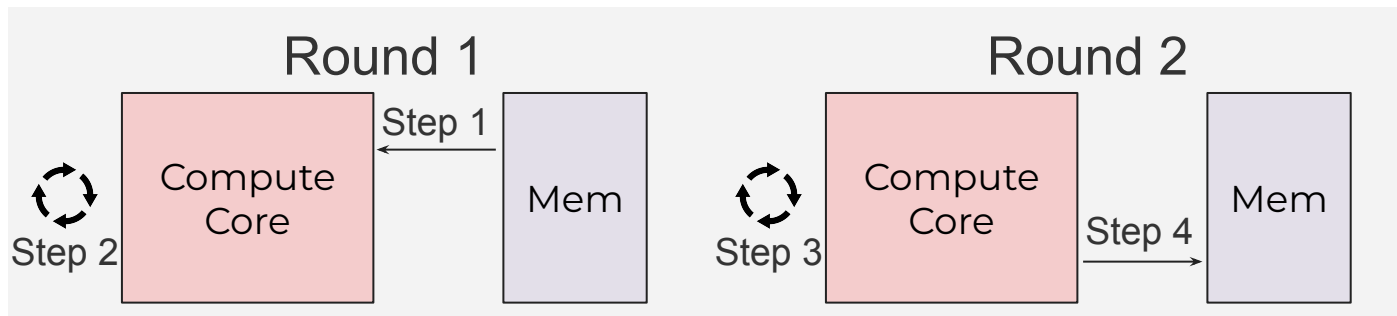
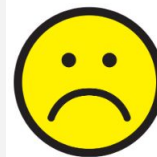
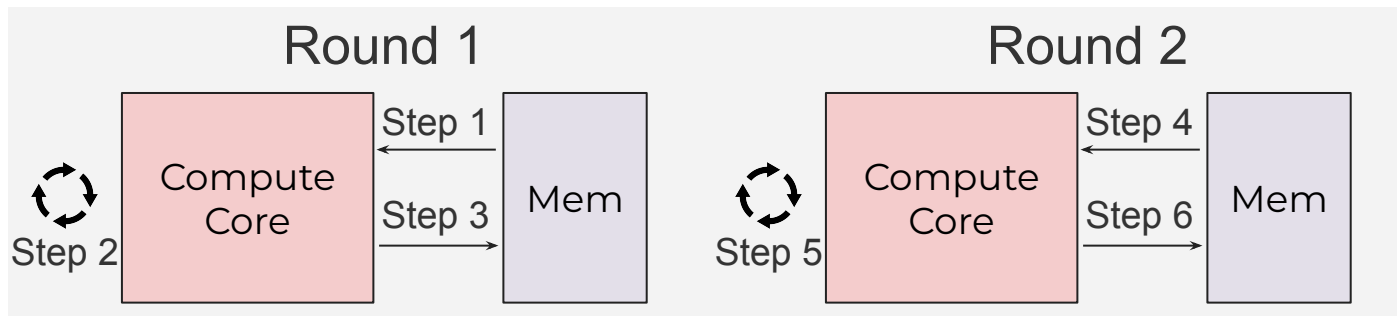
Flashattention



- Due to the heavy involvement of attention mechanism, transformers are memory-bound rather than compute bound.
- Kernel fusion: if there are multiple operations applied to the same input, the input can be loaded once from HBM, instead of multiple times for each operation.

Flashattention

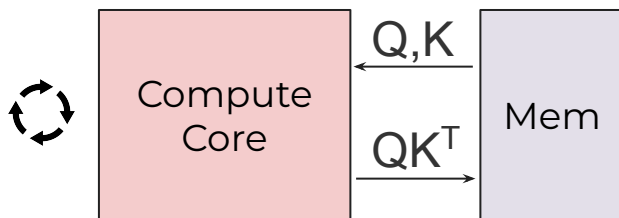
$$Y = S(QK^T) \times V$$



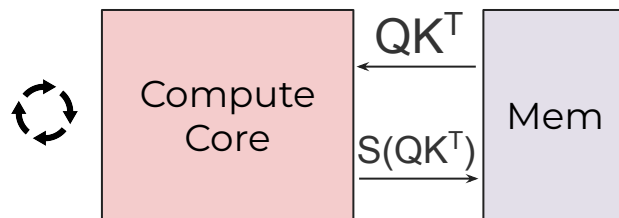
Flashattention

$$Y = S(QK^T) \times V$$

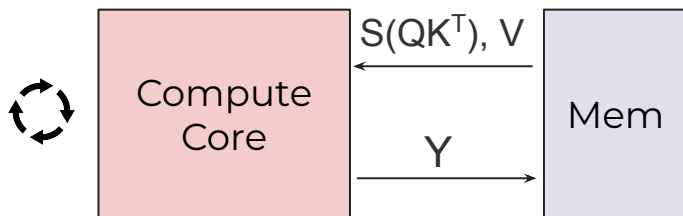
Round 1



Round 2



Round 3



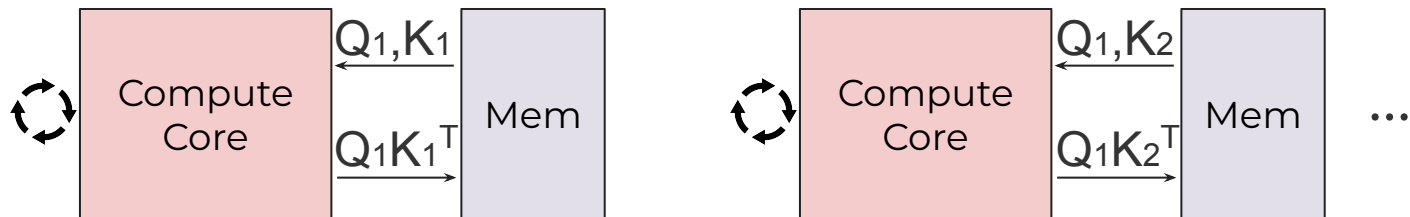
- The computation of QK^T must be all finished before computing softmax.
- This will lead to multiple rounds of memory access.

$$s_i = \frac{e^{z_i}}{\sum_{j=0}^{N-1} e^{z_j}} \text{ For } i = 1, 2, \dots, N$$

Flashattention

$$Y = S(QK^T) \times V$$

Round 1



- In practice, this operation will be executed in tiles.

Flashattention

Algorithm 0 Standard Attention Implementation

Require: Matrices $\mathbf{Q}, \mathbf{K}, \mathbf{V} \in \mathbb{R}^{N \times d}$ in HBM.

- 1: Load \mathbf{Q}, \mathbf{K} by blocks from HBM, compute $\mathbf{S} = \mathbf{Q}\mathbf{K}^\top$, write \mathbf{S} to HBM.
 - 2: Read \mathbf{S} from HBM, compute $\mathbf{P} = \text{softmax}(\mathbf{S})$, write \mathbf{P} to HBM.
 - 3: Load \mathbf{P} and \mathbf{V} by blocks from HBM, compute $\mathbf{O} = \mathbf{P}\mathbf{V}$, write \mathbf{O} to HBM.
 - 4: Return \mathbf{O} .
-

- Softmax and linear layers are computed separately.
- Flashattention splits the inputs $\mathbf{Q}, \mathbf{K}, \mathbf{V}$ into tiles, then compute the attention output with respect to those blocks.

Flashattention

- Softmax operation can be performed as:

$$m(x) := \max_i x_i, \quad f(x) := [e^{x_1 - m(x)} \quad \dots \quad e^{x_B - m(x)}]$$

$$\ell(x) := \sum_i f(x)_i, \quad \text{softmax}(x) := \frac{f(x)}{\ell(x)}$$

- We can fuse the softmax with the linear layer by doing follows:

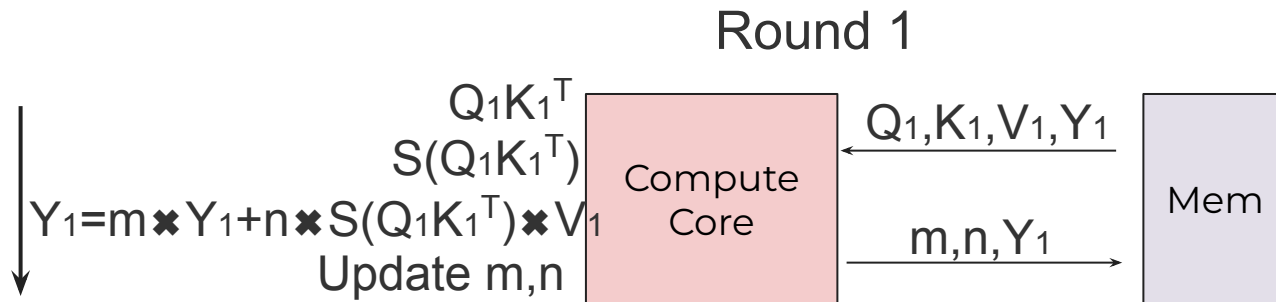
$$m(x) = m([x^{(1)} \quad x^{(2)}]) = \max(m(x^{(1)}), m(x^{(2)}))$$

$$f(x) = [e^{m(x^{(1)}) - m(x)} f(x^{(1)}) \quad e^{m(x^{(2)}) - m(x)} f(x^{(2)})]$$

$$\ell(x) = \ell([x^{(1)} \quad x^{(2)}]) = e^{m(x^{(1)}) - m(x)} \ell(x^{(1)}) + e^{m(x^{(2)}) - m(x)} \ell(x^{(2)}) \quad \text{softmax}(x) = \frac{f(x)}{\ell(x)}$$

Flashattention

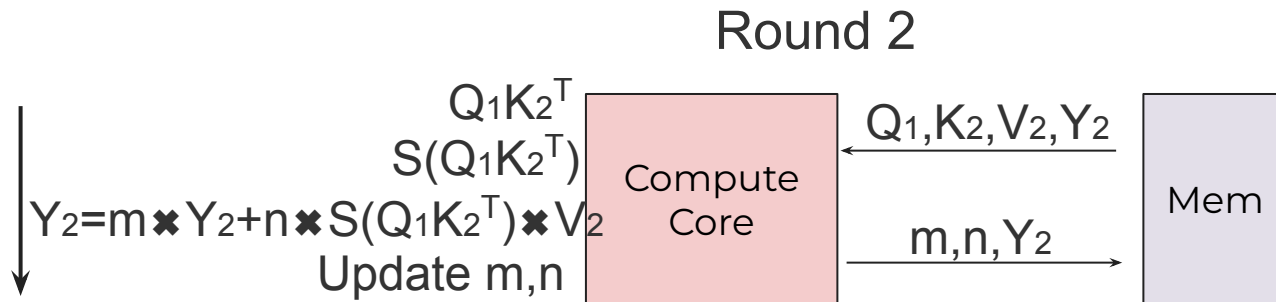
$$S(QK^T) \times V$$



$$s_i = \frac{e^{z_i}}{\sum_{j=0}^{N-1} e^{z_j}} \text{ For } i = 1, 2, \dots, N$$

Flashattention

$$S(QK^T) \times V$$



$$s_i = \frac{e^{z_i}}{\sum_{j=0}^{N-1} e^{z_j}} \text{ For } i = 1, 2, \dots, N$$

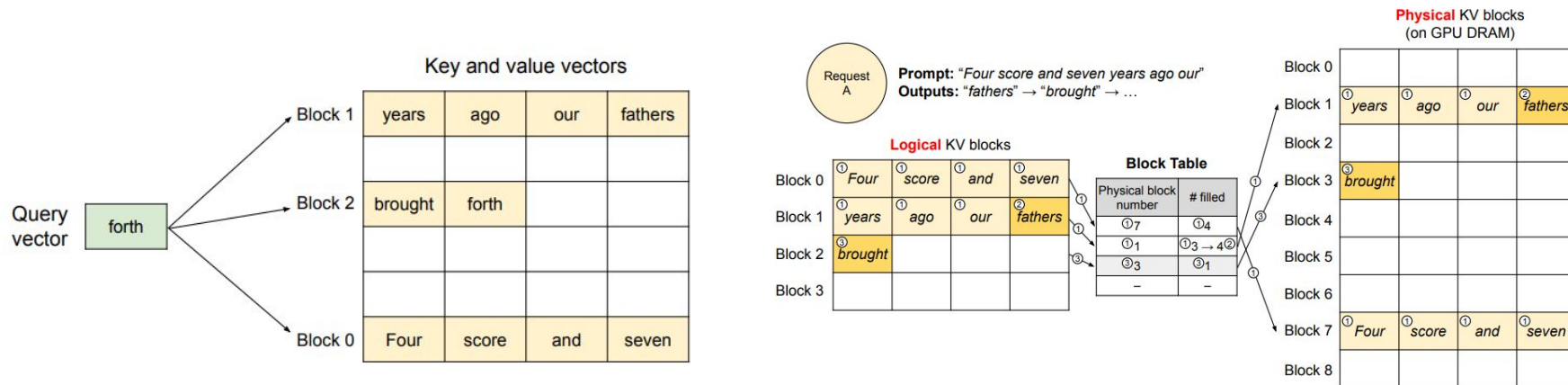
Flashattention

Algorithm 1 FLASHATTENTION

Require: Matrices $\mathbf{Q}, \mathbf{K}, \mathbf{V} \in \mathbb{R}^{N \times d}$ in HBM, on-chip SRAM of size M .

- 1: Set block sizes $B_c = \lceil \frac{M}{4d} \rceil$, $B_r = \min(\lceil \frac{M}{4d} \rceil, d)$.
 - 2: Initialize $\mathbf{O} = (0)_{N \times d} \in \mathbb{R}^{N \times d}$, $\ell = (0)_N \in \mathbb{R}^N$, $m = (-\infty)_N \in \mathbb{R}^N$ in HBM.
 - 3: Divide \mathbf{Q} into $T_r = \lceil \frac{N}{B_r} \rceil$ blocks $\mathbf{Q}_1, \dots, \mathbf{Q}_{T_r}$ of size $B_r \times d$ each, and divide \mathbf{K}, \mathbf{V} into $T_c = \lceil \frac{N}{B_c} \rceil$ blocks $\mathbf{K}_1, \dots, \mathbf{K}_{T_c}$ and $\mathbf{V}_1, \dots, \mathbf{V}_{T_c}$, of size $B_c \times d$ each.
 - 4: Divide \mathbf{O} into T_r blocks $\mathbf{O}_i, \dots, \mathbf{O}_{T_r}$ of size $B_r \times d$ each, divide ℓ into T_r blocks $\ell_i, \dots, \ell_{T_r}$ of size B_r each, divide m into T_r blocks m_1, \dots, m_{T_r} of size B_r each.
 - 5: **for** $1 \leq j \leq T_c$ **do**
 - 6: Load $\mathbf{K}_j, \mathbf{V}_j$ from HBM to on-chip SRAM.
 - 7: **for** $1 \leq i \leq T_r$ **do**
 - 8: Load $\mathbf{Q}_i, \mathbf{O}_i, \ell_i, m_i$ from HBM to on-chip SRAM.
 - 9: On chip, compute $\mathbf{S}_{ij} = \mathbf{Q}_i \mathbf{K}_j^T \in \mathbb{R}^{B_r \times B_c}$.
 - 10: On chip, compute $\tilde{m}_{ij} = \text{rowmax}(\mathbf{S}_{ij}) \in \mathbb{R}^{B_r}$, $\tilde{\mathbf{P}}_{ij} = \exp(\mathbf{S}_{ij} - \tilde{m}_{ij}) \in \mathbb{R}^{B_r \times B_c}$ (pointwise), $\tilde{\ell}_{ij} = \text{rowsum}(\tilde{\mathbf{P}}_{ij}) \in \mathbb{R}^{B_r}$.
 - 11: On chip, compute $m_i^{\text{new}} = \max(m_i, \tilde{m}_{ij}) \in \mathbb{R}^{B_r}$, $\ell_i^{\text{new}} = e^{m_i - m_i^{\text{new}}} \ell_i + e^{\tilde{m}_{ij} - m_i^{\text{new}}} \tilde{\ell}_{ij} \in \mathbb{R}^{B_r}$.
 - 12: Write $\mathbf{O}_i \leftarrow \text{diag}(\ell_i^{\text{new}})^{-1} (\text{diag}(\ell_i) e^{m_i - m_i^{\text{new}}} \mathbf{O}_i + e^{\tilde{m}_{ij} - m_i^{\text{new}}} \tilde{\mathbf{P}}_{ij} \mathbf{V}_j)$ to HBM.
 - 13: Write $\ell_i \leftarrow \ell_i^{\text{new}}$, $m_i \leftarrow m_i^{\text{new}}$ to HBM.
 - 14: **end for**
 - 15: **end for**
 - 16: Return \mathbf{O} .
-

Paged Attention

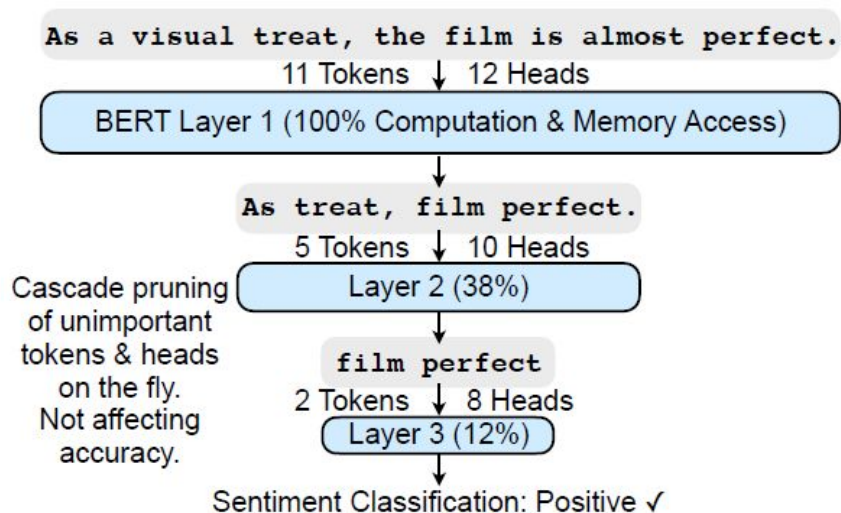


- An LLM serving system that achieves (1) near-zero waste in KV cache memory and (2) flexible sharing of KV cache within and across requests to further reduce memory usage.

Topics

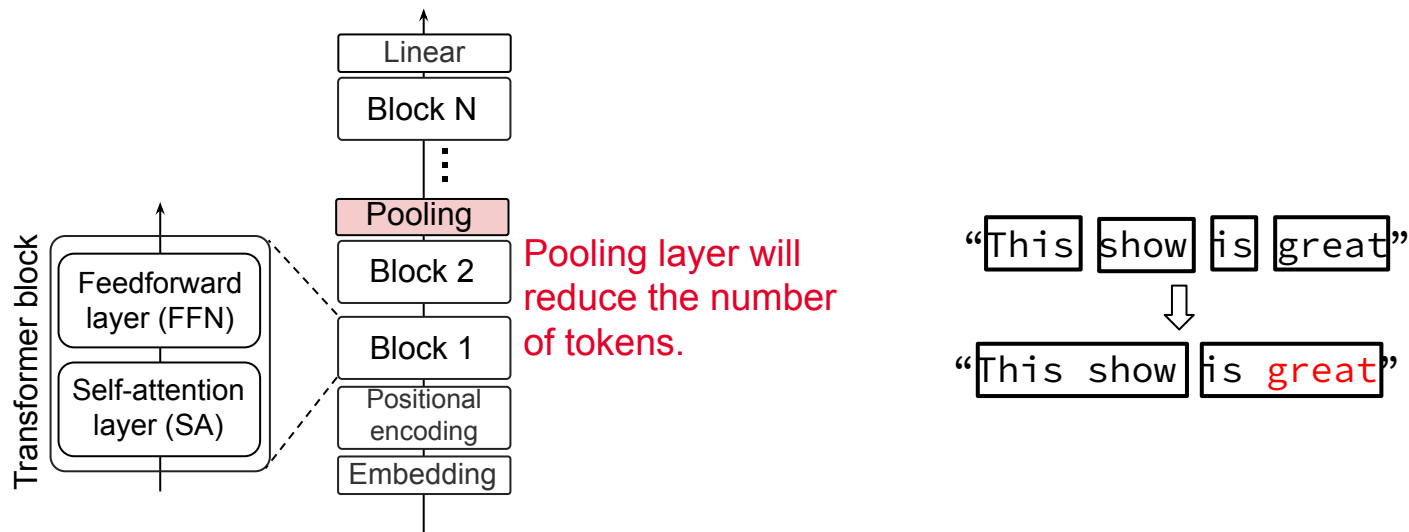
- Matrix Multiplication
- Hardware design for Nonlinear Blocks
- System optimization of LLMs
- Popular transformer accelerator design
 - SpAtten
 - EdgeBert
 - Olive

SpAtten: Efficient Sparse Attention Architecture with Cascade Token and Head Pruning



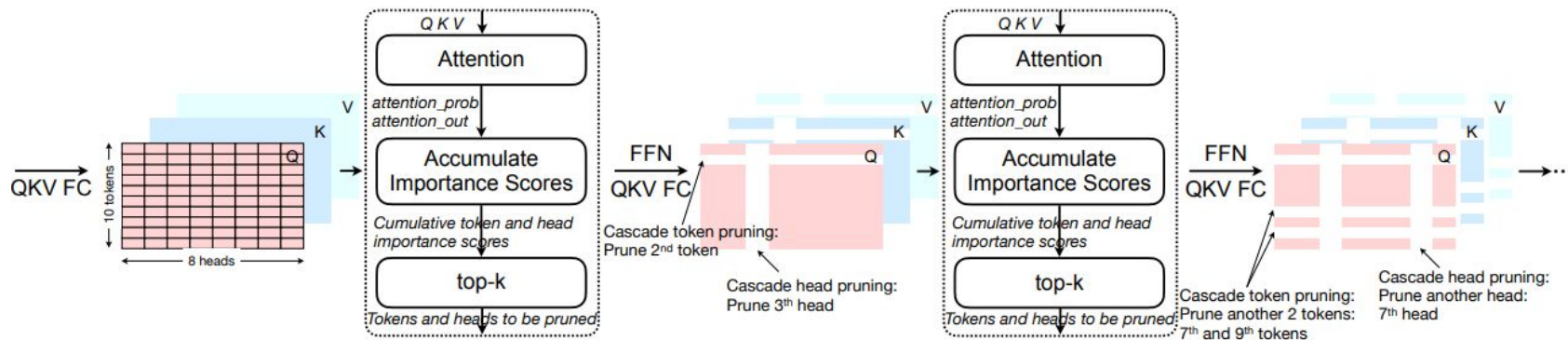
- Not all the tokens nor heads are necessary to produce the final results.
- SpAtten is an efficient algorithm-architecture co-design that leverages token sparsity, head sparsity, and quantization opportunities to reduce the attention computation and memory access.

Token Merging



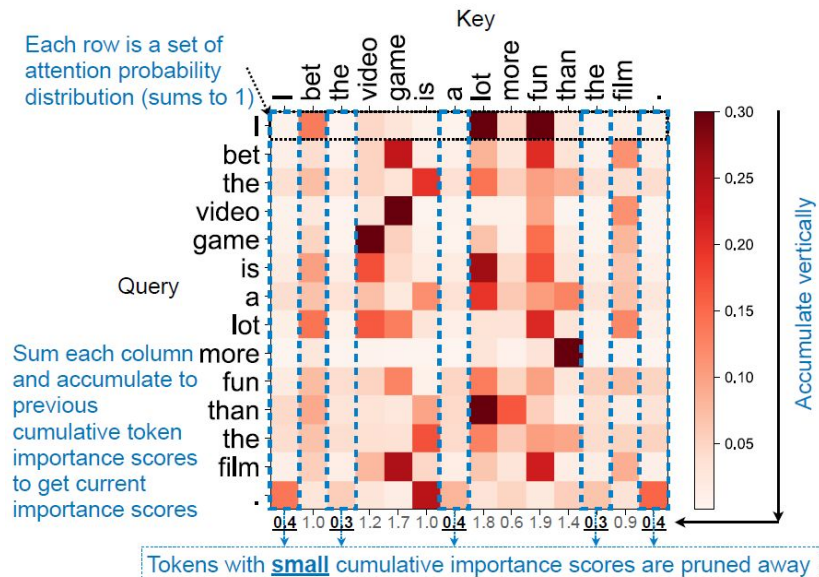
- We can reduce the number of tokens by merging them together.

SpAtten: Efficient Sparse Attention Architecture with Cascade Token and Head Pruning



- Cascade token pruning removes redundant tokens and corresponding entire Q K V vectors according to the cumulative token importance scores computed from attention prob.
- Cascade head pruning removes unimportant heads and corresponding chunks in all Q K V vectors according to the cumulative head important scores computed from attention out.

SpAtten: Efficient Sparse Attention Architecture with Cascade Token and Head Pruning



- For each input, the summation of each column indicates the importance of this token, we can remove the unimportant tokens accordingly.
- Similarly, we can compute the importance of each heads.

SpAtten: Efficient Sparse Attention Architecture with Cascade Token and Head Pruning

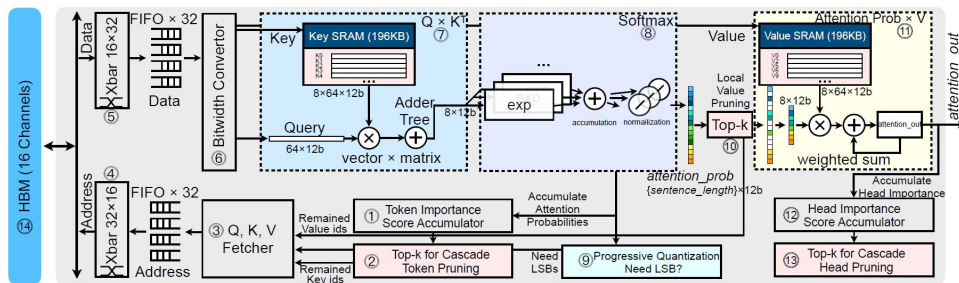
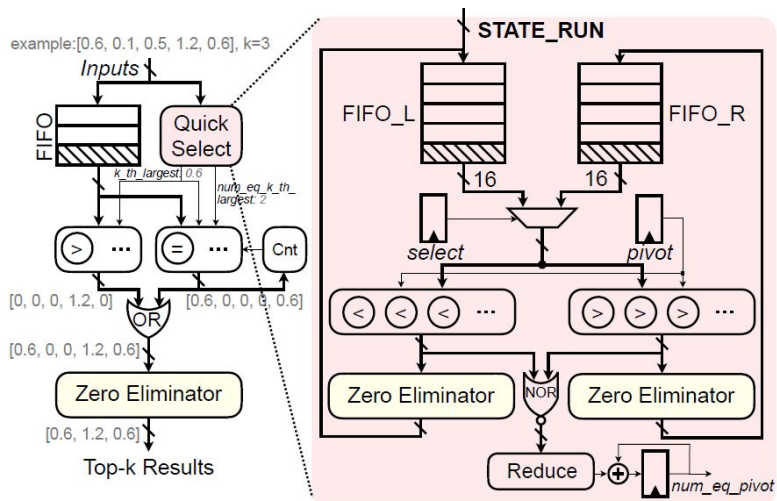
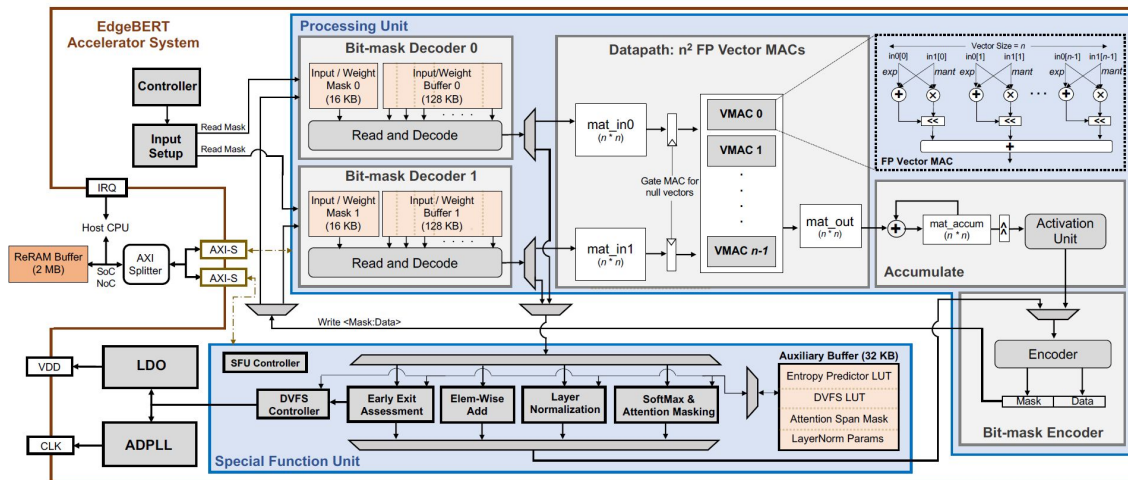


Fig. 8. SpAtten Architecture Overview. Modules on the critical path (6,7,8,10,11) are fully pipelined to maximize the throughput.

- This paper also proposed novel architecture to perform top-k extraction with high parallelism.



EdgeBert

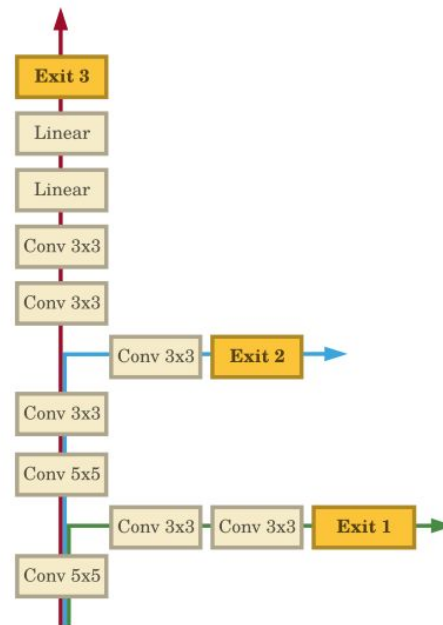


$$Energy \propto \alpha C V_{DD}^2 N_{cycles}$$

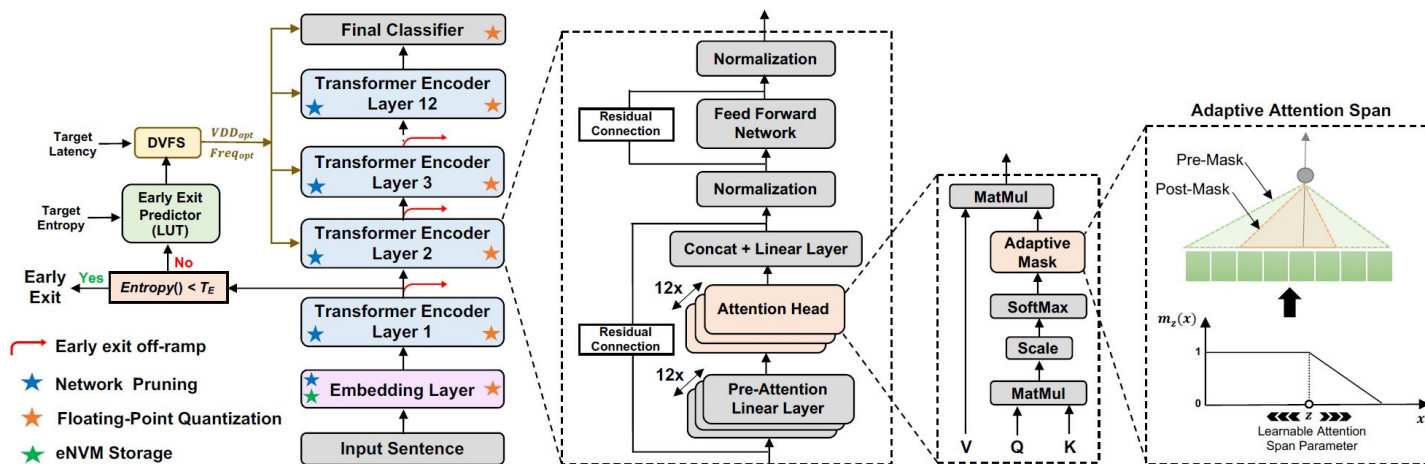
- EdgeBERT is a novel algorithm-hardware codesign approach to enable latency-bound NLP workloads on resource-constrained embedded devices.
- EdgeBERT dynamically tunes frequency and voltage settings to optimize the trade-off between accuracy, latency, and energy consumption.

Early Exit Mechanism

- EdgeBERT employs entropy-based early exit predication in order to perform dynamic voltage-frequency scaling (DVFS), at a sentence granularity, for minimal energy consumption while adhering to a prescribed target latency.
- During Inference, a confidence score is computed at each exit point, if greater than a predefined threshold, then the output is computed locally, leading to a faster inference.
- The confidence score is defined as: $\text{entropy}(\mathbf{y}) = \sum_{c \in \mathcal{C}} y_c \log y_c$,



Other Tricks for Efficiency



- Pruning and Quantization techniques are also adopted.

Latency Aware Inference

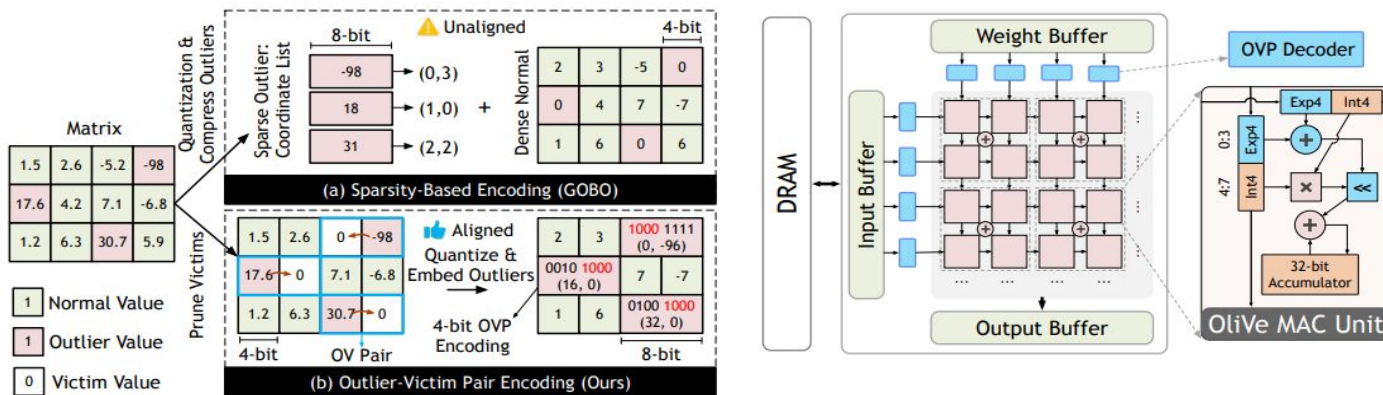
Algorithm 2: EdgeBERT latency-aware inference. Computations exit at the predicted exit layer or earlier.

Input: T := per-sentence latency target, E_T := entropy target

```
for input sentence  $i = 1$  to  $n$  do
  for encoder layer  $l = 1$  do
     $z_l = f(x; \theta | VDD_{nom}, Freq_{max})$ 
    if  $entropy(z_l) < E_T$  then
      exit inference
    else
       $L_{predict} = LUT(entropy(z_l), E_T)$ 
       $VDD_{opt}, Freq_{opt} = DVFS(L_{predict}, T)$ 
  for encoder layer  $l = 2$  to  $L_{predict}$  do
     $z_l = f(x; \theta | VDD_{opt}, Freq_{opt})$ 
    if  $entropy(z_l) < E_T$  then
      exit inference
  exit inference
```

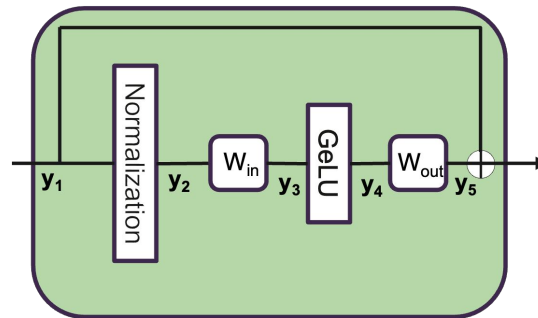
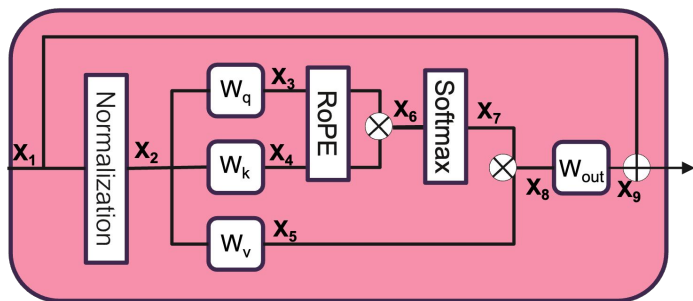
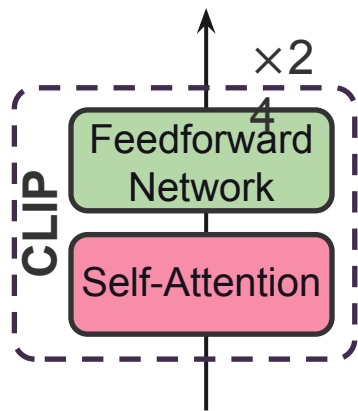
- The entropy result of the first layer is sent to a trained classifier (LUT-based) to predict which following encoder layer should early exit.
- The voltage and frequency is scaled down to proper energy-optimal setting for the rest of encoder layers while meeting the latency target for each sentence.
- This scheme produces a quadratic reduction in the accelerator power consumption.
- To realize fast per-sentence DVFS, the on-chip DVFS system is developed and integrated within EdgeBERT

OliVe: Accelerating LLMs via Hardware-friendly Outlier-Victim Pair Quantization



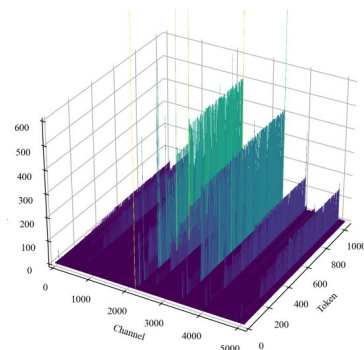
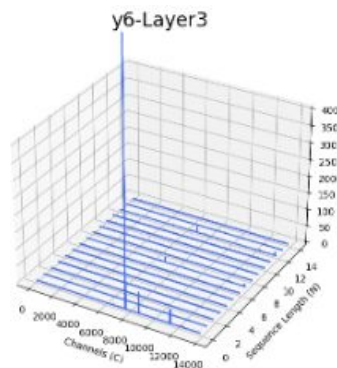
- Recent studies show when the model size exceeds a threshold (e.g., 6 billion), the model performance is vulnerable to only a tiny fraction ($< 0.1\%$) of outliers, whose values are much more significant than normal values.
- Olive adopts a hybrid quantization scheme and handles outlier values locally using a separate quantization scheme.

Outlier within LLMs

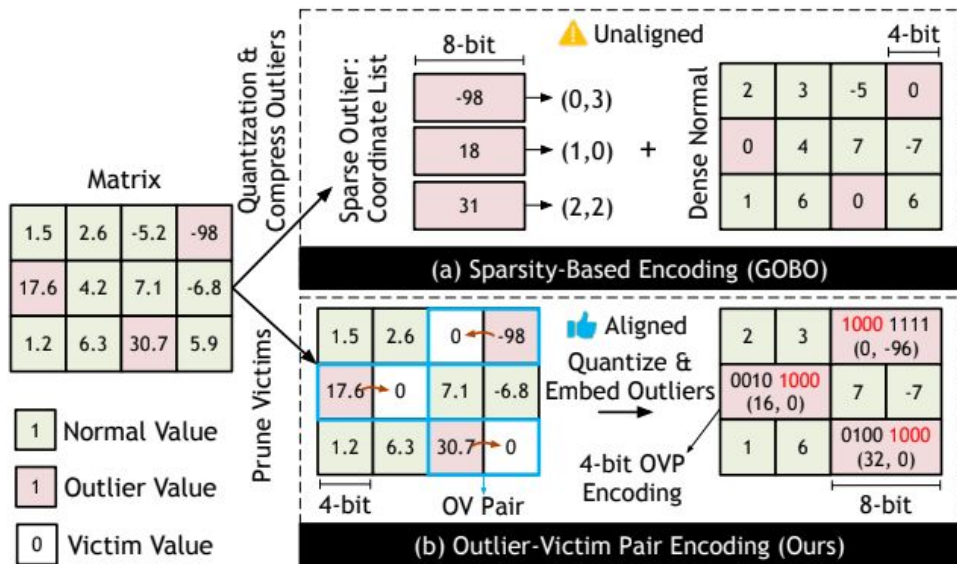


Types of Outlier

- Massive Activation:
 - For an activation matrix A , an massive activation is an element A_{ij} within it that satisfies:
 - $A_{ij} > \eta \times \text{mean}(|A|)$
 - $A_{ij} > \gamma$
 - $\eta=300, \gamma=50$
- Channelwise Outlier:
 - $\text{mean}(A_i) > \eta \times \text{std}(A) + \text{mean}(|A|)$
 - $\text{std}(A_i) < \beta$
 - $\eta=3, \beta=0.6$



OliVe: Accelerating LLMs via Hardware-friendly Outlier-Victim Pair Quantization



- In this work, we aim to design an architecture to handle outliers in a localized way with high hardware efficiency. Post quantization training is adopted.
- To better quantize the outlier, one number near the outlier are sacrificed and set to 0, then two set of bits are used to encode the outlier.

Presentations

- [Efficient Memory Management for Large Language Model Serving with PagedAttention](#) (Xiwen Min & Ziyun Cheng)
- [Cnvlutin: Ineffectual-Neuron-Free Deep Neural Network Computing](#) (Jishnu Warrior & Ishaan Shivhare)